

Linearly constrained Gaussian processes

Carl Jidling¹, Niklas Wahlström¹, Adrian Wills², and Thomas B. Schön¹

¹Department of Information Technology, Uppsala University, Sweden. E-mail: {carl.jidling, niklas.wahlstrom, thomas.schon}@it.uu.se

²School of Engineering, University of Newcastle, Australia. E-mail: adrian.wills@newcastle.edu.au

Abstract

We consider a modification of the covariance function in Gaussian processes to correctly account for known linear constraints. By modelling the target function as a transformation of an underlying function, the constraints are explicitly incorporated in the model such that they are guaranteed to be fulfilled by any sample drawn or prediction made. We also propose a constructive procedure for designing the transformation operator and illustrate the result on both simulated and real-data examples.

1 Introduction

Bayesian non-parametric modelling has had a profound impact in machine learning due, in no small part, to the flexibility of these model structures in combination with the ability to encode prior knowledge in a principled manner Ghahramani (2015). These properties have been exploited within the class of Bayesian non-parametric models known as Gaussian Processes (GPs), which have received significant research attention and have demonstrated utility across a very large range of real-world applications Rasmussen and Williams (2006).

Abstracting from the myriad number of these applications, it has been observed that the efficacy of GPs modelling is often intimately dependent on the appropriate choice of mean and covariance functions, and the appropriate tuning of their associated hyper-parameters. Often, the most appropriate mean and covariance functions are connected to prior knowledge of the underlying problem. For example, Koyejo et al. (2013) use functional expectation constraints to consider the problem of gene-disease association, and Navarro et al. (2016) employs a multivariate generalised von Mises distribution to produce a GP-like regression that handles circular variable problems.

At the same time, it is not always obvious how one might construct a GP model that obeys underlying principles, such as equilibrium conditions and conservation "laws". One straightforward approach to this problem is to add fictitious measurements that observe the constraints at a finite number of points of interest. This has the benefit of being relatively straightforward to implement, but has the sometimes significant drawback of increasing the problem dimension and at the same time not enforcing the constraints between the points of interest.

A different approach to constraining the GP model is to construct mean and covariance functions that obey the constraints. For example, curl and divergence free covariance functions are used by Wahlström (2015) to improve the accuracy for regression problems. The main benefit of this approach is that the problem dimension does not grow, and the constraints are enforced everywhere, not just at the points of interest. However, it is not obvious how these approaches can be easily scaled to handle an arbitrary set of linear operator constraints.

The contribution of this paper is a new way to include constraints into multivariate GPs. In particular, we develop an algorithm that transforms a given GP into a new, derived, one that satisfies linear operator constraints on the new GP. We will demonstrate the utility of this new method on both simulated examples and on a real-world application, the latter in form of predicting the components of a magnetic field, as illustrated in Figure 1.

To make these ideas more concrete, we present a simple example that will serve as a focal point several times throughout the paper. To that end, assume that we have a two-dimensional function $\mathbf{f}(\mathbf{x}) : \mathbb{R}^2 \mapsto \mathbb{R}^2$ on which we put a GP prior

$$\mathbf{f}(\mathbf{x}) \sim \mathcal{GP}(\boldsymbol{\mu}(\mathbf{x}), K(\mathbf{x}, \mathbf{x}')). \quad (1a)$$

We further know that $\mathbf{f}(\mathbf{x})$ should obey the differential equation

$$\frac{\partial f_1}{\partial x_1} + \frac{\partial f_2}{\partial x_2} = 0. \quad (1b)$$

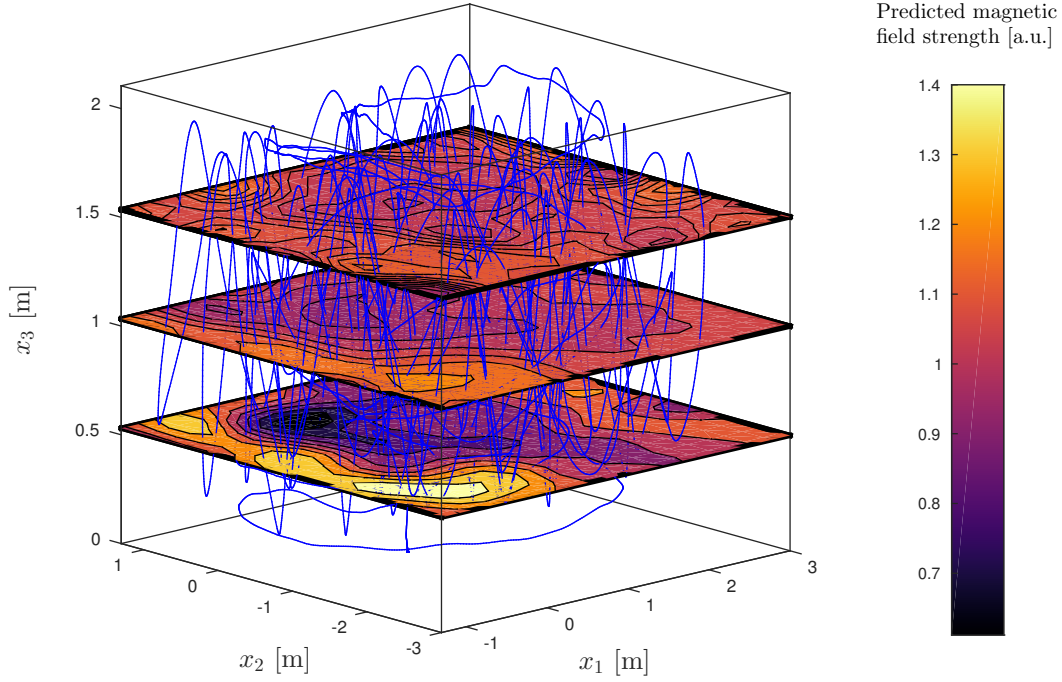


Figure 1: Predicted strength of a magnetic field at three heights, given measured data sampled from the trajectory shown (blue curve). The three components (x_1, x_2, x_3) denote the Cartesian coordinates, where the x_3 -coordinate is the height above the floor. The magnetic field is curl-free, which can be formulated in terms of three linear constraints. The method proposed in this paper can exploit these constraints to improve the predictions. See Section 5.2 for details.

The contribution of this paper is a procedure how to modify $K(\mathbf{x}, \mathbf{x}')$ and $\boldsymbol{\mu}(\mathbf{x})$ such that any sample from the new GP is guaranteed to obey the constraints like (1b), considering not just differentiation but any kind of linear operators.

2 Problem formulation

Assume that we are given a data set of N observations $\{\mathbf{x}_k, \mathbf{y}_k\}_{k=1}^N$ where \mathbf{x}_k denotes the input and \mathbf{y}_k the output. Both the input and output are potentially vector-valued, where $\mathbf{x}_k \in \mathbb{R}^D$ and $\mathbf{y}_k \in \mathbb{R}^K$.

We consider the regression problem where the data can be described by a non-parametric model

$$\mathbf{y}_k = \mathbf{f}(\mathbf{x}_k) + \mathbf{e}_k, \quad (2)$$

where \mathbf{e}_k is zero-mean white noise representing the measurement uncertainty. In this work, we place a vector-valued GP prior on \mathbf{f}

$$\mathbf{f}(\mathbf{x}) \sim \mathcal{GP}(\boldsymbol{\mu}(\mathbf{x}), K(\mathbf{x}, \mathbf{x}')), \quad (3)$$

with the mean function and the covariance function

$$\boldsymbol{\mu}(\cdot) : \mathbb{R}^D \mapsto \mathbb{R}^K, \quad (4a)$$

$$K(\cdot, \cdot) : \mathbb{R}^D \times \mathbb{R}^D \mapsto \mathbb{R}^K \times \mathbb{R}^K. \quad (4b)$$

Based on the data $\{\mathbf{x}_k, \mathbf{y}_k\}_{k=1}^N$, we would now like to find a posterior over the function $\mathbf{f}(\mathbf{x})$. In addition to the data, we know that the function \mathbf{f} should fulfill certain constraints

$$\mathcal{F}_{\mathbf{x}}[\mathbf{f}] = \mathbf{0}, \quad (5)$$

where \mathcal{F}_x is an operator mapping the function $\mathbf{f}(\mathbf{x})$ to another function $\mathbf{g}(\mathbf{x})$ as $\mathcal{F}_x[\mathbf{f}] = \mathbf{g}(\mathbf{x})$. We further require \mathcal{F}_x to be a linear operator meaning that

$$\mathcal{F}_x[\lambda_1 \mathbf{f}_1 + \lambda_2 \mathbf{f}_2] = \lambda_1 \mathcal{F}_x[\mathbf{f}_1] + \lambda_2 \mathcal{F}_x[\mathbf{f}_2], \quad (6)$$

where $\lambda_1, \lambda_2 \in \mathbb{R}$. The operator \mathcal{F} can for example be a linear transform $\mathcal{F}_x[\mathbf{f}] = C\mathbf{f}(\mathbf{x})$ which together with the constraint (5) forces a certain linear combination of the outputs to be linearly dependent.

The operator \mathcal{F}_x could also include other linear operations on the function $\mathbf{f}(\mathbf{x})$. For example, we might know that the function $\mathbf{f}(\mathbf{x}) : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ should obey a certain partial differential equation $\mathcal{F}_x[\mathbf{f}] = \frac{\partial f_1}{\partial x_1} + \frac{\partial f_2}{\partial x_2}$. A few more linear operators are listed in Appendix A, including integration as one the most well-known.

The constraints (5) can either come from known physical laws or other prior knowledge of the process generating the data. Our objective is to encode these constraints in the mean and covariance functions (4) such that any sample from the corresponding GP prior (3) always obeys the constraint (5).

3 Building a constrained Gaussian process

3.1 Approach based on artificial observations

Just as Gaussian distributions are closed under linear transformations, so are GPs closed under linear operations (see Appendix B). This can be used for a straightforward way of embedding linear operator constraints of the form (5) into GP regression. The idea is to treat the constraints as noise-free artificial observations $\{\tilde{\mathbf{x}}_k, \tilde{\mathbf{y}}_k\}_{k=1}^{\tilde{N}}$ with $\tilde{\mathbf{y}}_k = \mathbf{0}$ for all $k = 1 \dots \tilde{N}$. The regression is then performed on the model

$$\tilde{\mathbf{y}}_k = \mathcal{F}_{\tilde{\mathbf{x}}_k}[\mathbf{f}], \quad (7)$$

where $\tilde{\mathbf{x}}_k$ are input points in the domain of interest. For example, one could let these artificial input points $\tilde{\mathbf{x}}_k$ coincide with the points of prediction.

An advantage of this approach is that it allows constraints of the type (5) with a non-zero right hand side. Furthermore, there is no theoretical limit on how many constraints we can include (i.e. number of rows in \mathcal{F}_x) – although in practice, of course, there is.

However, this is problematic mainly for two reasons. First of all, it makes the problem size grow. Not only does this increase memory requirements and execution time, but it also results in a higher condition number and is thereby a potential source of numerical instability. This is especially clear from the fact that we want these observations to be noise-free, since the noise usually has a regularizing effect.

Secondly, the constraints are only enforced point-wise, and not continuously. Therefore, a sample drawn from the posterior will not fulfill the constraint anywhere but only in our chosen points. The obvious way of compensating for this is by increasing the number of points in which the constraints are observed – but that exacerbates the first problem. Clearly, the challenge grows quickly with the dimension of the inferred function.

Letting the constraints be a property of the covariance function removes these issues – it makes the enforcement continuous while the problem size is left unchanged. The question is, how can we design such a covariance function? This is addressed in the next section.

3.2 A new construction

We want to find a GP prior (3) such that any sample $\mathbf{f}(\mathbf{x})$ from that prior obeys the constraints (5). In turn, this leads to constraints on the mean and covariance functions (4) of that prior. However, instead of posing these constraints on the mean and covariance functions directly, we consider $\mathbf{f}(\mathbf{x})$ to be related to another function $\mathbf{g}(\mathbf{x})$ via some operator \mathcal{G}_x

$$\mathbf{f}(\mathbf{x}) = \mathcal{G}_x[\mathbf{g}]. \quad (8)$$

The constraints (5) then amounts to

$$\mathcal{F}_x[\mathcal{G}_x[\mathbf{g}]] = \mathbf{0}. \quad (9)$$

We would like this relation to be true for any function $\mathbf{g}(\mathbf{x})$. To do that, we will interpret \mathcal{F}_x and \mathcal{G}_x as matrices and use a similar procedure to that of solving systems of linear equations. Since \mathcal{F}_x and \mathcal{G}_x are linear operators, we can think of $\mathcal{F}_x[\mathbf{f}]$ and $\mathcal{G}_x[\mathbf{g}]$ as matrix-vector multiplications where

$$\mathcal{F}_x[\mathbf{f}] = \mathcal{F}_x \mathbf{f}, \quad \text{with} \quad (\mathcal{F}_x \mathbf{f})_i = \sum_{j=1}^K (\mathcal{F}_x)_{ij} f_j, \quad (10)$$

where each element $(\mathcal{F}_{\mathbf{x}})_{ij}$ in the operator matrix $\mathcal{F}_{\mathbf{x}}$ is a scalar operator. With this notation, (9) can be written as

$$\mathcal{F}_{\mathbf{x}} \mathcal{G}_{\mathbf{x}} = 0. \quad (11)$$

This reformulation imposes constraints on the operator $\mathcal{G}_{\mathbf{x}}$ rather than on the GP prior for $\mathbf{f}(\mathbf{x})$ directly.

We can now proceed by designing a GP prior for $\mathbf{g}(\mathbf{x})$ and transform it using the mapping (8). We further know that GPs are closed under linear operations. More specifically, if $\mathbf{g}(\mathbf{x})$ is modelled as a GP where

$$\mathbf{g}(\mathbf{x}) \sim \mathcal{GP}(\boldsymbol{\mu}_{\mathbf{g}}(\mathbf{x}), K_{\mathbf{g}}(\mathbf{x}, \mathbf{x}')), \quad (12)$$

then $\mathbf{f}(\mathbf{x})$ is also a GP with

$$\mathbf{f}(\mathbf{x}) = \mathcal{G}_{\mathbf{x}} \mathbf{g} \sim \mathcal{GP}(\mathcal{G}_{\mathbf{x}} \boldsymbol{\mu}_{\mathbf{g}}, \mathcal{G}_{\mathbf{x}} K_{\mathbf{g}} \mathcal{G}_{\mathbf{x}}^{\top}). \quad (13)$$

We use $(\mathcal{G}_{\mathbf{x}} K_{\mathbf{g}} \mathcal{G}_{\mathbf{x}'}^{\top})_{ij}$ to denote that

$$(\mathcal{G}_{\mathbf{x}} K_{\mathbf{g}} \mathcal{G}_{\mathbf{x}'}^{\top})_{ij} = (\mathcal{G}_{\mathbf{x}})_{ik} (\mathcal{G}_{\mathbf{x}'}^{\top})_{jl} (K_{\mathbf{g}})_{kl}, \quad (14)$$

where $\mathcal{G}_{\mathbf{x}}$ and $\mathcal{G}_{\mathbf{x}'}$ act on the first and second argument of $K_{\mathbf{g}}(\mathbf{x}, \mathbf{x}')$, respectively. See Appendix B for further details on linear operations on GPs.

The procedure to find the desired GP prior for \mathbf{f} can now be divided into the following three steps

1. Find an operator $\mathcal{G}_{\mathbf{x}}$ that fulfills the condition (9).
2. Choose a mean and covariance function for $\mathbf{g}(\mathbf{x})$.
3. Find the mean and covariance functions for $\mathbf{f}(\mathbf{x})$ according to (13).

In addition to being resistant to the disadvantages of the approach described in Section 3.1, there are some additional strengths worth pointing out with this method.

First of all, we have separated the task of encoding the constraints and encoding other desired properties of the kernel. The GP prior for $\mathbf{f}(\mathbf{x})$ will inherit the properties of the prior for $\mathbf{g}(\mathbf{x})$, such as those embedded in smoothness assumptions. In other words, we do not need to sacrifice any desired behaviour of the target function in order for the constraints to be fulfilled.

Secondly, $K(\mathbf{x}, \mathbf{x}')$ is guaranteed to be a valid covariance function, provided that $K_{\mathbf{g}}(\mathbf{x}, \mathbf{x}')$ is. This is due to the fact that GPs are closed under linear functional transformations. From (13), it is clear that each column of K must fulfill all constraints encoded in $\mathcal{F}_{\mathbf{x}}$. It may be possible to construct K only based on this knowledge, assuming a general form and solving the resulting system of equations. However, a solution may not just be hard to find, but one must also make sure that it is indeed a valid covariance function.

Furthermore, this approach provides a simple and straight-forward way of constructing the covariance function even if the constraints have a complicated form. It makes no difference if the linear operators relate the components of the target function explicitly or implicitly – the procedure remains the same.

3.3 Illustrating example

We will now illustrate the method using the example (1) introduced already in the introduction. Consider a function $\mathbf{f}(\mathbf{x}) : \mathbb{R}^2 \mapsto \mathbb{R}^2$ satisfying

$$\frac{\partial f_1}{\partial x_1} + \frac{\partial f_2}{\partial x_2} = 0, \quad (15)$$

where $\mathbf{x} = [x_1, x_2]^{\top}$ and $\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), f_2(\mathbf{x})]^{\top}$. This equation describes all two-dimensional divergence-free vector fields. Divergence-free vector fields arise, for example, in electromagnetics and fluid dynamics. The constraint (15) can be written as a linear constraint on the form (5) where

$$\underbrace{\begin{bmatrix} \frac{\partial}{\partial x_1} & \frac{\partial}{\partial x_2} \end{bmatrix}}_{\mathcal{F}_{\mathbf{x}}} \underbrace{\begin{bmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \end{bmatrix}}_{\mathbf{f}(\mathbf{x})} = \mathcal{F}_{\mathbf{x}} \mathbf{f} = 0.$$

Modelling this function with a GP and building the covariance structure as described above, we first need to find the transformation $\mathcal{G}_{\mathbf{x}}$ such that (11) is fulfilled. For example, we could pick

$$\mathcal{G}_{\mathbf{x}} = \begin{bmatrix} -\frac{\partial}{\partial x_2} & \frac{\partial}{\partial x_1} \end{bmatrix}^{\top}. \quad (16)$$

If the underlying function $g(\mathbf{x}) : \mathbb{R}^2 \mapsto \mathbb{R}$ is given by

$$g(\mathbf{x}) \sim \mathcal{GP}(0, k_g(\mathbf{x}, \mathbf{x}')),$$

we can make use of (13) to obtain

$$\mathbf{f}(\mathbf{x}) \sim \mathcal{GP}(\mathbf{0}, K(\mathbf{x}, \mathbf{x}')),$$

where

$$\begin{aligned} K(\mathbf{x}, \mathbf{x}') &= \mathcal{G}_{\mathbf{x}} k_g(\mathbf{x}, \mathbf{x}') \mathcal{G}_{\mathbf{x}}^\top \\ &= \begin{bmatrix} \frac{\partial^2}{\partial x_2 x_2'} & -\frac{\partial^2}{\partial x_2 x_1'} \\ -\frac{\partial^2}{\partial x_1 x_2'} & \frac{\partial^2}{\partial x_1 x_1'} \end{bmatrix} k_g(\mathbf{x}, \mathbf{x}'). \end{aligned}$$

Using a covariance function with the following structure, we know that the constraint (15) will be fulfilled by any function generated from the corresponding GP.

4 Finding the operator $\mathcal{G}_{\mathbf{x}}$

In a general setting it might be hard to find an operator $\mathcal{G}_{\mathbf{x}}$ that fulfills the constraint (11). Ultimately, we want an algorithm that can construct $\mathcal{G}_{\mathbf{x}}$ from a given $\mathcal{F}_{\mathbf{x}}$.

In more formal terms, the function $\mathcal{G}_{\mathbf{x}} \mathbf{g}$ forms the nullspace of $\mathcal{F}_{\mathbf{x}}$. The concept of nullspaces for linear operators is well-established Luenberger (1969), and does in many ways relate to real-number linear algebra.

However, an important difference can be illustrated by considering a one-dimensional function $f(x)$ subject to the constraint $\mathcal{F}_x f = 0$ where $\mathcal{F}_x = \frac{\partial}{\partial x}$. The solution to this differential equation can not be expressed in terms of an arbitrary underlying function, but we know that the relation is true if $f(x) = C$ where C is a scalar constant. Hence, the nullspace of $\frac{\partial}{\partial x}$ consists of the set of horizontal lines. Compare this with the real number equation $ab = 0$, $a \neq 0$, which is true if and only if $b = 0$. Since the nullspace differs between operators, it is clear that we must be careful when discussing the properties of $\mathcal{F}_{\mathbf{x}}$ and $\mathcal{G}_{\mathbf{x}}$ based on our knowledge from real-number algebra.

Let us denote the rows in $\mathcal{F}_{\mathbf{x}}$ as $\mathbf{f}_1^\top, \dots, \mathbf{f}_L^\top$. We now want to find all solutions \mathbf{g} such that

$$\mathcal{F}_{\mathbf{x}} \mathbf{g} = \mathbf{0} \quad \Rightarrow \quad \mathbf{f}_i^\top \mathbf{g} = 0, \quad \forall \quad i = 1, \dots, L. \quad (17)$$

The solutions $\mathbf{g}_1, \dots, \mathbf{g}_P$ to (17) will then be the columns of $\mathcal{G}_{\mathbf{x}}$. Each row vector \mathbf{f}_j can be written as $\mathbf{f}_j = \Phi_i \boldsymbol{\xi}^f$ where $\Phi_i \in \mathbb{R}^{K \times M_f}$ and $\boldsymbol{\xi}^f = [\xi_1, \dots, \xi_{M_f}]^\top$ is a vector of M_f scalar operators included in $\mathcal{F}_{\mathbf{x}}$.

We now assume that \mathbf{g} also can be written in a similar form $\mathbf{g} = \Gamma \boldsymbol{\xi}^g$ where $\Gamma \in \mathbb{R}^{K \times M_g}$ and $\boldsymbol{\xi}^g = [\xi_1, \dots, \xi_{M_g}]^\top$ is a vector of M_g scalar operators. One may make the assumption that the same set of operators that are used to describe \mathbf{f}_i also can be used to describe \mathbf{g} , i.e., $\boldsymbol{\xi}^g = \boldsymbol{\xi}^f$. However, this assumption might need to be relaxed.

The constraints (17) can then be written as

$$(\boldsymbol{\xi}^f)^\top \Phi_i \Gamma \boldsymbol{\xi}^g = 0, \quad \forall \quad i = 1, \dots, L. \quad (18)$$

We perform the multiplication and collect the terms in $\boldsymbol{\xi}^f$ and $\boldsymbol{\xi}^g$. The condition (18) then results in conditions on the parameters in Γ resulting in a homogeneous system of linear equations

$$A \cdot \text{vec}(\Gamma) = \mathbf{0}. \quad (19)$$

The vectors $\text{vec}(\Gamma_1), \dots, \text{vec}(\Gamma_P)$ spanning the nullspace of A in (19) are then used to compute the columns in $\mathcal{G}_{\mathbf{x}} = [\mathbf{g}_1, \dots, \mathbf{g}_P]$ where $\mathbf{g}_p = \Gamma_p \boldsymbol{\xi}^g$. If it turns out that the nullspace of A is empty, one should start over with a new ansatz and extend the set of operators in $\boldsymbol{\xi}^g$.

The outline of the procedure as described above is summarised in Algorithm 1. Let us now illustrate the method with an example.

The algorithm for finding $\mathcal{G}_{\mathbf{x}}$ is based upon a parametric ansatz rather than directly upon the theory for linear operators. Not only is it more intuitive, but it does also remove any conceptual challenges that theory may provide. A problem with this is that one may have to iterate before having found the appropriate set of operators in \mathcal{G} . It might be of interest to examine possible alternatives to this algorithm that does not use a parametric approach.

Algorithm 1 Constructing \mathcal{G}_x

Input: Operator matrix \mathcal{F}_x

Output: Operator matrix \mathcal{G}_x where $\mathcal{F}_x \mathcal{G}_x = \mathbf{0}$

Step 1: Make an ansatz $\mathbf{g} = \Gamma \xi^{\mathbf{g}}$ for the columns in \mathcal{G}_x .

Step 2: Expand $\mathcal{F}_x \Gamma \xi^{\mathbf{g}}$ and collect terms.

Step 3: Set up the system of equations $A \cdot \text{vec}(\Gamma) = \mathbf{0}$ and find the vectors $\Gamma_1 \dots \Gamma_P$ spanning its nullspace.

Step 4: If $P = 0$, go back to **Step 1** and make a new ansatz, i.e. extend the set of operators.

Step 5: Construct $\mathcal{G}_x = [\Gamma_1 \xi^{\mathbf{g}}, \dots, \Gamma_P \xi^{\mathbf{g}}]$.

4.1 Divergence-free example revisited

Let us return to the example discussed in Section 3.3, and show how the solution found by visual inspection also can be found with the algorithm described above.

Since \mathcal{F}_x only contains first-order derivative operators, we assume that a column in \mathcal{G}_x does so as well. Hence, let us propose the following ansatz (step 1)

$$\mathbf{g} = \begin{bmatrix} \gamma_{11} & \gamma_{12} \\ \gamma_{21} & \gamma_{22} \end{bmatrix} \begin{bmatrix} \frac{\partial}{\partial x_1} \\ \frac{\partial}{\partial x_2} \end{bmatrix} = \Gamma \xi^{\mathbf{g}}. \quad (20)$$

Applying the constraint, expanding and collecting terms (step 2) we find

$$\begin{aligned} \mathcal{F}_x \Gamma \xi^{\mathbf{g}} &= \begin{bmatrix} \frac{\partial}{\partial x_1} & \frac{\partial}{\partial x_2} \end{bmatrix} \begin{bmatrix} \gamma_{11} & \gamma_{12} \\ \gamma_{21} & \gamma_{22} \end{bmatrix} \begin{bmatrix} \frac{\partial}{\partial x_1} \\ \frac{\partial}{\partial x_2} \end{bmatrix} \\ &= \gamma_{11} \frac{\partial^2}{\partial x_1^2} + \gamma_{12} \frac{\partial^2}{\partial x_1 \partial x_2} + \gamma_{21} \frac{\partial^2}{\partial x_2 \partial x_1} + \gamma_{22} \frac{\partial^2}{\partial x_2^2} \\ &= \gamma_{11} \frac{\partial^2}{\partial x_1^2} + (\gamma_{12} + \gamma_{21}) \frac{\partial^2}{\partial x_1 \partial x_2} + \gamma_{22} \frac{\partial^2}{\partial x_2^2}, \end{aligned} \quad (21)$$

where we have used the fact that $\frac{\partial^2}{\partial x_i \partial x_j} = \frac{\partial^2}{\partial x_j \partial x_i}$ assuming continuous second derivatives. The expression (21) equals zero if

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \gamma_{11} \\ \gamma_{12} \\ \gamma_{21} \\ \gamma_{22} \end{bmatrix} = A \cdot \text{vec}(\Gamma) = \mathbf{0}. \quad (22)$$

The nullspace is spanned by a single vector (step 3)

$$\begin{bmatrix} \gamma_{11} \\ \gamma_{12} \\ \gamma_{21} \\ \gamma_{22} \end{bmatrix} = \lambda \begin{bmatrix} 0 \\ -1 \\ 1 \\ 0 \end{bmatrix}, \quad \lambda \in \mathbb{R}. \quad (23)$$

Choosing $\lambda = 1$, we get (step 5)

$$\mathcal{G}_x = \begin{bmatrix} -\frac{\partial}{\partial x_2} & \frac{\partial}{\partial x_1} \end{bmatrix}^T, \quad (24)$$

which is the same as in (16).

4.2 Generalisation

Although there are no conceptual problems with the algorithm introduced above, the procedure of expanding and collecting terms appears a bit informal. In a general form, the algorithm is reformulated such that the operators are completely left out from the solution process. The drawback of this is a more cumbersome notation, and we have therefore limited the presentation to this simplified version. However, the general algorithm is found in the Appendix C.

5 Experimental results

5.1 Simulated divergence-free function

Consider the example described in Section 3.3. An example of a function fulfilling the constraint (15) is

$$\begin{aligned} f_1(x_1, x_2) &= e^{-ax_1x_2} (ax_1 \sin(x_1x_2) - x_1 \cos(x_1x_2)), \\ f_2(x_1, x_2) &= e^{-ax_1x_2} (x_2 \cos(x_1x_2) - ax_2 \sin(x_1x_2)), \end{aligned} \quad (25)$$

where a denotes a constant.

We will now study how the regression of this function differs when using the covariance function found in Section 3.3 as compared to a diagonal covariance function

$$K(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}, \mathbf{x}')I.$$

The measurements generated are corrupted with Gaussian noise

$$\mathbf{y}_k = \mathbf{f}(\mathbf{x}_k) + \mathbf{e}_k, \quad \mathbf{e}_k \sim \mathcal{N}(\mathbf{0}, \sigma^2 I). \quad (26)$$

The squared exponential covariance function $k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp[-\frac{1}{2}l^{-2}\|\mathbf{x} - \mathbf{x}'\|^2]$ has been used for k_g and k with hyperparameters chosen by maximising the marginal likelihood. We have used the value $a = 0.01$ in (25).

We have included the approach described in Section 3.1 and varied the number of artificial observations to see the effect of this with respect to the error as compared to the other approaches.

We have used 50 measurements randomly picked over the domain $[0, 4] \times [0, 4]$, generated with the noise level $\sigma = 10^{-4}$. The points for prediction corresponds to a discretization using 20 uniformly distributed points in each direction, and hence a total of $N_P = 20^2 = 400$. The artificial observations have been chosen as random subsets of the prediction points, up to and including the full set.

The comparison is made with regard to the root mean squared error

$$e_{\text{rms}} = \sqrt{\frac{1}{N_P} \bar{\mathbf{f}}_{\Delta}^T \bar{\mathbf{f}}_{\Delta}}, \quad (27)$$

where $\bar{\mathbf{f}}_{\Delta} = \hat{\hat{\mathbf{f}}} - \bar{\mathbf{f}}$ and

$$\bar{\mathbf{f}} = [f_1(\mathbf{x}_1) \quad f_2(\mathbf{x}_1) \quad f_1(\mathbf{x}_2) \quad f_2(\mathbf{x}_2) \quad \dots]^T \quad (28)$$

is a concatenated vector storing the true function values in all prediction points and $\hat{\hat{\mathbf{f}}}$ denotes the reconstructed equivalent. To decrease the impact of randomness, each error value has been formed as an average over 50 reconstructions given different sets of measurements.

An example of the true field, measured values and reconstruction errors using the different methods is seen in Figure 2. The result from the experiment is seen in Figure 3. Note that the error from the approach with artificial observations is decreasing as the number of observations is increased, but only to a certain point. Have in mind, however, that the Gram matrix is growing, making the problem larger and worse conditioned. The result from our approach is clearly better, while the problem size is kept small and numerical problems are therefore avoided.

5.2 Real data experiment

Magnetic fields can mathematically be considered as a vector field mapping a 3D position to a 3D magnetic field strength. Based on the magnetostatic equations, this can be modeled as a curl-free vector field. By following Appendix C.1, our method can be used to encode the constraints in the following covariance function

$$K_{\text{curl}}(\mathbf{x}, \mathbf{x}') = \sigma_f^2 e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2l^2}} \left(I_3 - \left(\frac{\mathbf{x} - \mathbf{x}'}{l} \right) \left(\frac{\mathbf{x} - \mathbf{x}'}{l} \right)^T \right), \quad (29)$$

which also has been presented elsewhere Wahlström (2015).

With a magnetic sensor and an optical positioning system, both position and magnetic field data have been collected in a magnetically distorted indoor environment, see Appendix D for details about the experimental details. This data was collected together with a data set previously published by Solin et al. (2015). In Figure 1 the predicted magnitude of the magnetic field over a two-dimensional domain for three different heights above the floor is displayed. The predictions have been made based on 500 measurements sampled from the trajectory that is given by the blue curve.

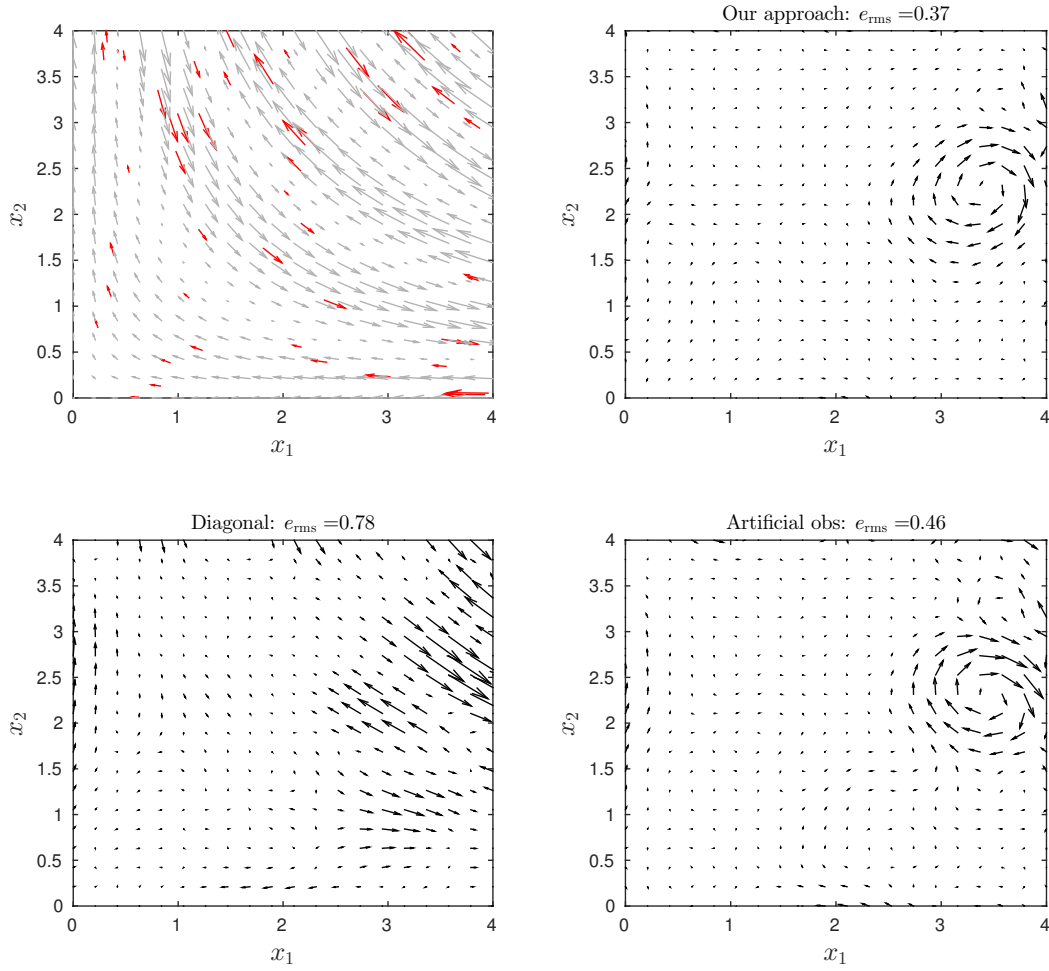


Figure 2: Top left: Example of field plots illustrating the measurements (red arrows) and the true field (grey arrows). Remaining three plots: reconstructed fields subtracted from the true field over the domain of prediction. The artificial observations of the constraint have been made in the same points as the predictions are made.

Similar to the simulated experiment in Section 5.1, we compare the predictions of the curl-free covariance function (29) with the diagonal covariance function and the diagonal covariance function using artificial observations. The results have been formed by averaging the error over 50 reconstructions. In each iteration, training data and test data were randomly selected from the data set collected in the experiment. 500 train data points and 1 000 test data points were used.

The result is seen in Figure 4. We recognise the same behaviour as we saw for the simulated experiment in Figure 3. Note that the accuracy of the artificial observation approach gets very close to our approach for a large number of artificial observations. However, in the last step of increasing the artificial observations, the accuracy decreases. This is probably caused by the numerical errors that follows from an ill-conditioned Gram matrix.

6 Related work

Many problems in which GPs are used contain some kind of constraint that could be well exploited to improve the quality of the solution. Since there are a variety of ways in which constraints may appear and take form, there is also a variety of methods to deal with them.

The treatment of inequality constraints in GP regression have been considered for instance by Abrahamsen and

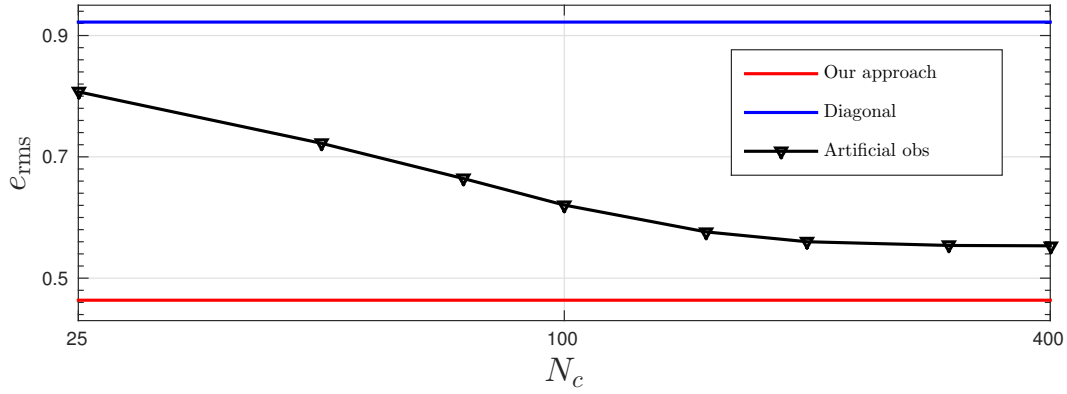


Figure 3: Accuracy of the different approaches as the number of artificial observations N_c is increased for the simulated experiment.

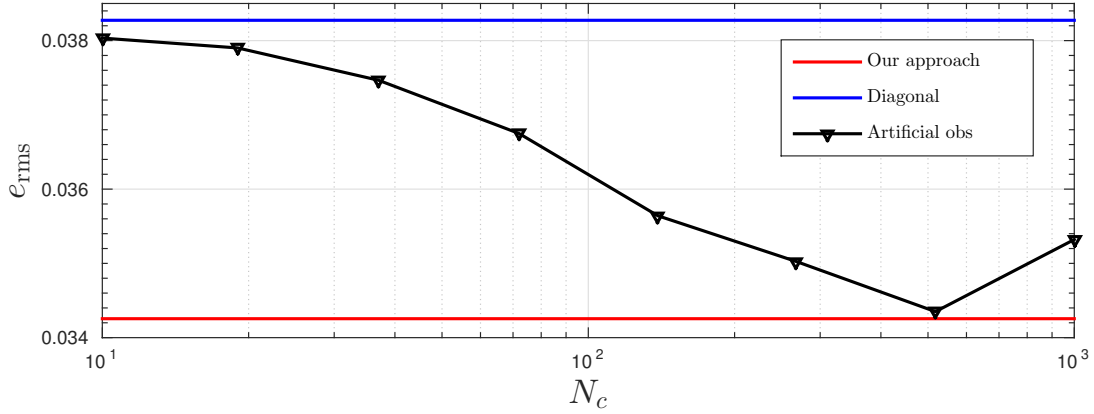


Figure 4: Accuracy of the different approaches as the number of artificial observations N_c is increased for the real-data experiment.

Benth (2001) and Da Veiga and Marrel (2012), based on local representation in a limited set of points. The paper by Maatouk and Bay (2016) proposes a finite-dimensional GP-approximation to allow for inequality constraints in the entire domain.

It has been shown that linear constraints satisfied by the training data will be satisfied by the GP prediction as well Salzmann and Urtasun (2010). The same paper shows how this result can be extended to quadratic forms through a parametric reformulation and minimisation of the Frobenious norm, with application demonstrated for pose estimation. Another approach on capturing human body features is described by Rudovic and Pantic (2011), where a face-shape model is included in the GP framework to imply anatomic correctness.

Although constraints in most situations are formulated on the outputs of the GP, there are also situations in which they are acting on the inputs. An example of this is given by Tran et al. (2015), describing a method of benefit from ordering constraints on the input to reduce the negative impact of input noise.

Within medicine, examples of applications include gene-disease association through functional expectation constraints Koyejo et al. (2013) and lung disease sub-type identification using a mixture of GPs and constraints encoded with Markov random fields Ross and Dy (2013).

Another way of viewing constraints is as modified prior distributions. For example, by making use of the so-called multivariate generalised von Mises distribution, Navarro et al. (2016) ends up in a version of GP regression customised for circular variable problems.

Generally speaking, the papers mentioned above consider problems in which the constraints are dealt with using

some kind of external enforcement – that is, they are not explicitly incorporated into the model, but rely on approximations or finite representations. Therefore, the constraints may just be approximately satisfied and not necessarily in a continuous manner, which differs from the method proposed in this paper. Of course, comparisons can not be done directly between methods that have been developed for different kinds of constraints. The interest in this paper is multivariate problems where the constraints are linear combinations of the outputs that are known to equal zero.

For multivariate problems, the construction of the covariance function does include an extra challenge in modelling the correlation between the output components. Different techniques are described in the literature, and we refer to Álvarez et al. (2012) for a useful review. The basic idea behind the so-called *separable kernels* is to separate the process of modelling the covariance function for each component and the process of modelling the correlation between them. The final covariance function can then be chosen for example according to some method of regularisation.

Another class of covariance functions is referred to as the *invariant kernels*. Here, the correlation is inherited from a known mathematical relation. The curl- and divergence free covariance functions are such examples where the structure follows directly from the underlying physics, and has been shown to improve the accuracy notably for regression problems Wahlström (2015). Another example is the method proposed by Constantinescu and Anitescu (2013), where the Taylor expansion is used to construct a covariance model given a known relationship between the outputs. Other fields of interest include using GPs in solving partial differential equations Graepel (2003); Nguyen and Peraire (2015), as well as other kinds of functional regression Nguyen and Peraire (2016).

A very useful property on linear transformations is given by Särkkä (2011), based on the GPs natural inheritance of features imposed by linear operators. This fact has for example been used in developing a method for monitoring infectious diseases Andrade-Pacheco et al. (2016).

The method proposed in this work is exploiting the transformation property to build a covariance function of the invariant kind. Unlike aforementioned work, the cornerstone here is the choice of the transformation, which is made such that the constraints are built into the prior and therefore are guaranteed to be fulfilled.

7 Conclusion and future work

In this article we have presented a method for designing the covariance function of a multivariate Gaussian process subject to known linear operator constraints on the target function. The method will by construction guarantee that any sample drawn from the resulting process will obey the constraints in all points. Numerical simulations show the benefits of this method as compared to alternative approaches. Furthermore, it has been demonstrated to improve the performance on real data as well.

As mentioned in Section 4, it would be desirable to describe the requirements on \mathcal{G}_x more rigorously. If that is possible it might allow us to reformulate the algorithm for the construction of \mathcal{G}_x in a way that allows for a more straightforward approach as compared to the parametric ansatz that we have proposed. In particular, our method relies upon the requirement that the target function can be expressed in terms of an underlying *potential* function g . This leads to the intriguing and nontrivial question: Is it possible to mathematically guarantee the existence of such a potential? If the answer to this question is yes, the next question will of course be what it look like and what is its relationship to the target function.

Another possible topic of further research is the extension to constraints including nonlinear operators, which for example might rely upon a linearisation in the domain of interest. Furthermore, it may be of potential interest to study the case of a non-zero right-hand side of (5), which clearly would require some modifications.

8 Acknowledgements

This research is financially supported by the Swedish Foundation for Strategic Research (SSF) via the project *ASSEMBLE* (Contract number: RIT 15-0012). The work is also supported by the Swedish Research Council (VR) via the project *Probabilistic modeling of dynamical systems* (Contract number: 621-2013-5524). We are grateful for the help and equipment provided by the UAS Technologies Lab, Artificial Intelligence and Integrated Computer Systems Division (AIICS) at the Department of Computer and Information Science (IDA), Linköping University, Sweden. The real data set used in this paper has been collected by some of the authors together with Manon Kok, Arno Solin, and Simo Särkkä. We thank them for allowing us to use this data. We also thank Manon Kok for supporting us with the data processing. Furthermore, we would like to thank Carl Rasmussen and Marc Deisenroth for fruitful discussions on constrained GPs.

A Linear operators

In this work we consider linear operators on functions. Such an operator transforms a function $\mathbf{f}(\mathbf{x})$ to another function $\mathbf{g}(\mathbf{z})$. We denote this according to

$$\mathbf{g}(\mathbf{z}) = \mathcal{F}_{\mathbf{z}}[\mathbf{f}(\mathbf{x})]. \quad (30)$$

This linear operator could be *differentiation* of a function. If $D = 1$ and $K = 1$ this will be defined as

$$g(z) = \mathcal{F}_z[f] = \left. \frac{\partial f(x)}{\partial x} \right|_{x=z} \quad (31a)$$

which slightly more informal also can be written as

$$g(x) = \mathcal{F}_x[f] = \frac{\partial f(x)}{\partial x}. \quad (31b)$$

Also *integration* of a scalar function $f(x)$ over an interval $[z_1, z_2]$ is a linear operator

$$g(\mathbf{z}) = \mathcal{F}_{\mathbf{z}}[f] = \int_{z_1}^{z_2} f(x) dx, \quad (32)$$

where $g(\mathbf{z})$ is a scalar-valued function with a two-dimensional input $\mathbf{z} = [z_1, z_2]^\top$. Note that in the two examples given above, the inputs of f and g will not be the same, not even of the same dimension!

Input wrapping is another way to construct new covariance functions from old ones (Rasmussen and Williams, 2006, page 92). It utilizes a nonlinear wrapping $\mathbf{x} = \mathbf{u}(\mathbf{z})$ of the input variables. This wrapping can also be considered as a linear operator, where

$$\mathbf{g}(\mathbf{z}) = \mathcal{F}_{\mathbf{z}}[\mathbf{f}] = \mathbf{f}(\mathbf{x})|_{\mathbf{x}=\mathbf{u}(\mathbf{z})}. \quad (33)$$

This operator also changes the function input and possibly also its dimension. Even though the wrapping itself might be nonlinear, the operator corresponding to this wrapping is in fact linear.

It is straightforward to show that all three operators presented above do fulfill the linearity condition (6).

B Gaussian processes under linear operations

It is well-known that Gaussian distributions are closed under linear transformation. In similar manner, Gaussian processes are closed under linear operations (Papoulis and Pillai, 1991; Rasmussen and Williams, 2006; Hennig and Kiefel, 2013; Garnett, 2015).

By applying the functional $\mathcal{F}_{\mathbf{x}}$ on both the mean function and the covariance function, the GP prior for $\mathcal{F}_{\mathbf{x}}$ is given by

$$\mathcal{F}_{\mathbf{x}}\mathbf{f} \sim \mathcal{GP}(\mathcal{F}_{\mathbf{x}}\boldsymbol{\mu}, \text{Cov}[\mathcal{F}_{\mathbf{x}}\mathbf{f}(\mathbf{x}), \mathcal{F}_{\mathbf{x}'}\mathbf{f}(\mathbf{x}')]). \quad (34)$$

The covariance becomes

$$\begin{aligned} & \text{Cov}[\mathcal{F}_{\mathbf{x}}\mathbf{f}(\mathbf{x}), \mathcal{F}_{\mathbf{x}'}\mathbf{f}(\mathbf{x}')] \\ &= \mathbb{E} \left[(\mathcal{F}_{\mathbf{x}}\mathbf{f}(\mathbf{x}) - \mathcal{F}_{\mathbf{x}}\boldsymbol{\mu}(\mathbf{x})) (\mathcal{F}_{\mathbf{x}'}\mathbf{f}(\mathbf{x}') - \mathcal{F}_{\mathbf{x}'}\boldsymbol{\mu}(\mathbf{x}'))^\top \right] \\ &= \mathcal{F}_{\mathbf{x}} \mathbb{E} \left[(\mathbf{f}(\mathbf{x}) - \boldsymbol{\mu}(\mathbf{x})) (\mathbf{f}(\mathbf{x}') - \boldsymbol{\mu}(\mathbf{x}'))^\top \right] \mathcal{F}_{\mathbf{x}'}^\top \\ &= \mathcal{F}_{\mathbf{x}} K \mathcal{F}_{\mathbf{x}'}^\top, \end{aligned} \quad (35)$$

where by the notation $(\mathcal{F}_{\mathbf{x}} K \mathcal{F}_{\mathbf{x}'}^\top)_{ij}$ we mean that

$$(\mathcal{F}_{\mathbf{x}} K \mathcal{F}_{\mathbf{x}'}^\top)_{ij} = (\mathcal{F}_{\mathbf{x}})_{ik} (\mathcal{F}_{\mathbf{x}'}^\top)_{jl} K_{kl}, \quad (36)$$

and where $(\mathcal{F}_{\mathbf{x}})_{ik}$ and $(\mathcal{F}_{\mathbf{x}'}^\top)_{jl}$ act on the first and second argument of $K_{kl}(\mathbf{x}, \mathbf{x}')$, respectively.

We should point out that some care must be taken when applying this procedure. For example, if we would like to consider the derivative of a function governed by a GP, we must make sure that this function is modelled in a way such that the derivative actually exists. This may sound obvious, yet important to remember since the set of standard covariance functions includes members that are not differentiable – among those we find Matérn_{1/2} Rasmussen and Williams (2006).

C Generalization of Section 4

In this supplementary material we will generalize the method described in the main paper on how to solve operator matrix equations on the form

$$\mathcal{F}\mathcal{G} = \mathbf{0},$$

where we want to find \mathcal{G} given \mathcal{F} ¹. If $\mathcal{F} \in \mathbb{R}^{m \times n}$ is a real valued matrix, \mathcal{G} can easily be found by letting the columns in \mathcal{G} span the nullspace of \mathcal{F} (provided such a nullspace exist). However, if the elements of \mathcal{F} are operators, the situation is more tricky. This supplementary material generalizes the parametric approach presented in Section 4 in the main paper for arbitrary operators of any order. The strategy is to study the vector space of homogeneous polynomials where the operators are interpreted as the variables of these polynomials.

In Section C.1, we assume that both \mathcal{F} and \mathcal{G} consist of first order operators and in Section C.2 we generalize this to allow for any order of the operators.

C.1 First order operator equation

Consider the matrix $\mathcal{F} \in \mathcal{P}_p^{m \times n}$, where \mathcal{P}_p is a vector space of first order operators

$$\mathcal{P}_p = \{a_1 y_1 + \dots + a_p y_p \mid a_1, \dots, a_p \in \mathbb{R}\}, \quad (37)$$

where y_1, \dots, y_p is the basis in that vector space. The basis components y_k can for example represent derivative operators $y_k = \frac{\partial}{\partial x_k}$. We want to find the vectors $\mathbf{g} \in \mathcal{P}_p^n$ such that $\mathcal{F}\mathbf{g} = \mathbf{0}$ is fulfilled. We can write $\mathcal{F} \in \mathcal{P}_p^{m \times n}$ and $\mathbf{g} \in \mathcal{P}_p^n$ as

$$\mathcal{F}_{ij} = \sum_{k=1}^p \phi_{ijk} y_k, \quad \phi_{ijk} = \{\Phi\}_{ijk} \in \mathbb{R}, \quad (38a)$$

$$\mathbf{g}_j = \sum_{k=1}^p \gamma_{jk} y_k, \quad \gamma_{jk} = \{\Gamma\}_{jk} \in \mathbb{R}, \quad (38b)$$

where $\Phi \in \mathbb{R}^{m \times n \times p}$ and $\Gamma \in \mathbb{R}^{n \times p}$. This gives

$$\mathcal{F}\mathbf{g} = \mathbf{0} \Leftrightarrow \sum_{j=1}^n \sum_{k=1}^p \sum_{l=1}^p \phi_{ijk} y_k \gamma_{jl} y_l = 0 \quad \forall i = 1 : m. \quad (39)$$

For each i , we have a quadratic form

$$\mathbf{y}^\top \Phi_i \Gamma \mathbf{y} = 0, \quad (40)$$

where $\Phi_i \in \mathbb{R}^{p \times n}$ with $\{\Phi_i\}_{kj} = \phi_{ijk}$ and $\Gamma \in \mathbb{R}^{n \times p}$ with $\{\Gamma\}_{jk} = \gamma_{jk}$.

The quadratic form is equal to zero for all \mathbf{y} if and only if

$$\Phi_i \Gamma + \Gamma^\top \Phi_i^\top = 0 \quad \forall i = 1 : m. \quad (41)$$

Example 1 (divergence free vector field)

We consider the following vector of operators $\mathcal{F} \in \mathcal{P}_3^{1 \times 3}$

$$\mathcal{F} = \nabla_{\mathbf{x}} = \left[\frac{\partial}{\partial x_1}, \frac{\partial}{\partial x_2}, \frac{\partial}{\partial x_3} \right], \quad (42)$$

where

$$\mathcal{F}_{ij} = \sum_{k=1}^3 \phi_{ijk} y_k, \quad \forall i = 1, \quad j = 1, 2, 3, \quad (43)$$

¹In this appendix, the argument \mathbf{x} is omitted for simplified notation

where $y_k = \frac{\partial}{\partial x_k}$. Following the notation introduced above, for this particular operator matrix we have

$$\Phi_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}. \quad (44)$$

We now want to find a vector $\mathbf{g} \in \mathcal{P}^3$ that fulfills $\mathcal{F}\mathbf{g} = \mathbf{0}$ for all \mathbf{y} . We assume that this operator vector is in \mathcal{P}_3^3 and can be written

$$\mathbf{g}_j = \sum_{k=1}^3 \gamma_{jk} y_k \quad j = 1, 2, 3, \quad (45)$$

where $\Gamma \in \mathbb{R}^{3 \times 3}$ is unknown. Now we have that

$$\Phi_1 \Gamma + \Gamma^\top \Phi_1^\top = 0 \quad (46a)$$

$$\Rightarrow \begin{bmatrix} \gamma_{11} & \gamma_{12} - \gamma_{21} & \gamma_{13} - \gamma_{31} \\ \gamma_{21} - \gamma_{12} & \gamma_{22} & \gamma_{23} - \gamma_{32} \\ \gamma_{31} - \gamma_{13} & \gamma_{32} - \gamma_{23} & \gamma_{33} \end{bmatrix} = 0, \quad (46b)$$

which in turn gives

$$\gamma_{11} = 0, \quad \gamma_{12} + \gamma_{21} = 0, \quad (47a)$$

$$\gamma_{22} = 0, \quad \gamma_{13} + \gamma_{31} = 0, \quad (47b)$$

$$\gamma_{33} = 0, \quad \gamma_{23} + \gamma_{32} = 0. \quad (47c)$$

The nullspace of (46a) is then spanned by

$$\Gamma = \lambda_1 \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & -1 & 0 \end{bmatrix} + \lambda_2 \begin{bmatrix} 0 & 0 & -1 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix} + \lambda_3 \begin{bmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix},$$

which gives

$$\mathbf{g} = \lambda_1 \begin{bmatrix} 0 \\ \frac{\partial}{\partial x_3} \\ -\frac{\partial}{\partial x_2} \end{bmatrix} + \lambda_2 \begin{bmatrix} -\frac{\partial}{\partial x_3} \\ 0 \\ \frac{\partial}{\partial x_1} \end{bmatrix} + \lambda_3 \begin{bmatrix} \frac{\partial}{\partial x_2} \\ -\frac{\partial}{\partial x_1} \\ 0 \end{bmatrix}, \lambda_1, \lambda_2, \lambda_3 \in \mathbb{R}.$$

Example 2 (curl free vector field)

We consider the following vector of operators $\mathcal{F} \in \mathcal{P}_3^{3 \times 3}$

$$\mathcal{F} = \begin{bmatrix} 0 & \frac{\partial}{\partial x_3} & -\frac{\partial}{\partial x_2} \\ -\frac{\partial}{\partial x_3} & 0 & \frac{\partial}{\partial x_1} \\ \frac{\partial}{\partial x_2} & -\frac{\partial}{\partial x_1} & 0 \end{bmatrix}, \quad (48)$$

where

$$\mathcal{F}_{ij} = \sum_{k=1}^3 \phi_{ijk} y_k, \quad \forall i = 1:3, \quad j = 1:3, \quad (49)$$

where $y_k = \frac{\partial}{\partial x_k}$. For this particular operator matrix we have

$$\Phi_1 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 1 & 0 \end{bmatrix}, \quad \Phi_2 = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ -1 & 0 & 0 \end{bmatrix}, \quad \Phi_3 = \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

We now want to find a vector $\mathbf{g} \in \mathcal{P}^3$ which fulfills $\mathcal{F}\mathbf{g} = \mathbf{0}$ for all \mathbf{y} . We assume that this operator vector is in \mathcal{P}_3^3 and can be written

$$\mathbf{g}_j = \sum_{k=1}^3 \gamma_{jk} y_k \quad j = 1, 2, 3, \quad (50)$$

where $\Gamma \in \mathbb{R}^{3 \times 3}$ is unknown. Now we have that

$$\begin{aligned}\Phi_1 \Gamma + \Gamma^\top \Phi_1^\top &= 0 \Rightarrow \begin{bmatrix} 0 & -\gamma_{31} & \gamma_{21} \\ -\gamma_{31} & -2\gamma_{32} & \gamma_{22}-\gamma_{33} \\ \gamma_{21} & \gamma_{22}-\gamma_{33} & 2\gamma_{23} \end{bmatrix} = 0, \\ \Phi_2 \Gamma + \Gamma^\top \Phi_2^\top &= 0 \Rightarrow \begin{bmatrix} 2\gamma_{31} & \gamma_{32} & \gamma_{33}-\gamma_{11} \\ \gamma_{32} & 0 & -\gamma_{12} \\ \gamma_{33}-\gamma_{11} & -\gamma_{12} & -2\gamma_{13} \end{bmatrix} = 0, \\ \Phi_3 \Gamma + \Gamma^\top \Phi_3^\top &= 0 \Rightarrow \begin{bmatrix} 2\gamma_{21} & \gamma_{22}-\gamma_{11} & \gamma_{23} \\ \gamma_{22}-\gamma_{11} & -2\gamma_{12} & -\gamma_{13} \\ \gamma_{23} & -\gamma_{13} & 0 \end{bmatrix} = 0,\end{aligned}$$

which in turn gives

$$\gamma_{22} - \gamma_{33} = 0, \quad \gamma_{23} = 0, \quad \gamma_{32} = 0, \quad (51a)$$

$$\gamma_{33} - \gamma_{11} = 0, \quad \gamma_{13} = 0, \quad \gamma_{31} = 0, \quad (51b)$$

$$\gamma_{22} - \gamma_{11} = 0, \quad \gamma_{12} = 0, \quad \gamma_{21} = 0. \quad (51c)$$

The nullspace of (51a) is then spanned by the single base vector

$$\Gamma = \lambda_1 \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad \lambda_1 \in \mathbb{R}, \quad (52)$$

which gives

$$\mathbf{g} = \lambda_1 \begin{bmatrix} \frac{\partial}{\partial x_1} \\ \frac{\partial}{\partial x_2} \\ \frac{\partial}{\partial x_3} \end{bmatrix}, \quad \lambda_1 \in \mathbb{R}. \quad (53)$$

The final covariance function becomes

$$K(\mathbf{x}, \mathbf{x}') = \begin{bmatrix} \frac{\partial^2}{\partial x_1 \partial x'_1} & \frac{\partial^2}{\partial x_1 \partial x'_2} & \frac{\partial^2}{\partial x_1 \partial x'_3} \\ \frac{\partial^2}{\partial x_2 \partial x'_1} & \frac{\partial^2}{\partial x_2 \partial x'_2} & \frac{\partial^2}{\partial x_2 \partial x'_3} \\ \frac{\partial^2}{\partial x_3 \partial x'_1} & \frac{\partial^2}{\partial x_3 \partial x'_2} & \frac{\partial^2}{\partial x_3 \partial x'_3} \end{bmatrix} k_g(\mathbf{x}, \mathbf{x}'). \quad (54)$$

If we use the squared exponential covariance function

$$k_g(\mathbf{x}, \mathbf{x}') = \sigma_f^2 e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2l^2}} \quad (55)$$

we get

$$K(\mathbf{x}, \mathbf{x}') = \frac{\sigma_f^2}{l^2} e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2l^2}} \left(I_3 - \left(\frac{\mathbf{x} - \mathbf{x}'}{l} \right) \left(\frac{\mathbf{x} - \mathbf{x}'}{l} \right)^\top \right). \quad (56)$$

This covariance function is used in the real data experiment in Section 5.2 of the main paper. Note, that the version in the paper does not use l^2 in the denominator (which we also would get here if we would multiply (53) with l^2 , still providing the same constraints).

C.2 Higher order operator equation

Now, consider the matrix $\mathcal{F} \in \mathcal{P}_{p,q}^{m \times n}$, where $\mathcal{P}_{p,q}$ is a vector space of all homogeneous polynomials of degree q in p variables

$$\mathcal{P}_{p,q} = \left\{ \sum_{k_1}^p \cdots \sum_{k_q}^p a_{k_1, \dots, k_q} y_{k_1} \cdots y_{k_q} \Gamma i g | a_{k_1, \dots, k_q} \in \mathbb{R} \right\},$$

where the nominals $y_{k_1} \cdots y_{k_q}$ constitute the basis of that vector space. The components y_k can for example represent derivative operators $y_k = \frac{\partial}{\partial x_k}$ and $\mathcal{P}_{p,q}$ then contain all q th order derivatives of $x_1 \dots x_q$. We want to find the vectors $\mathbf{g} \in \mathcal{P}_{p,q_g}^n$ such that $\mathcal{F}\mathbf{g} = \mathbf{0}$ is fulfilled. We can write $\mathcal{F} \in \mathcal{P}_{p,q}^{m \times n}$ and $\mathbf{g} \in \mathcal{P}_{p,q_g}^n$ as

$$\mathcal{F}_{ij} = \sum_{k_1}^p \cdots \sum_{k_q}^p \phi_{i,j,k_1,\dots,k_q} y_{k_1} \cdots y_{k_q}, \quad (57a)$$

$$\mathbf{g}_j = \sum_{k_1}^p \cdots \sum_{k_q}^p \gamma_{j,k_1,\dots,k_q} y_{k_1} \cdots y_{k_q}, \quad (57b)$$

where $\Phi \in \mathbb{R}^{m \times n \times p^{\times q}}$ and $\mathbf{b} \in \mathbb{R}^{n \times p^{\times q}}$ (here $p^{\times q}$ denotes $\underbrace{p \times \cdots \times p}_{q \text{ times}}$). This gives

$$\mathcal{F}\mathbf{g} = 0 \Leftrightarrow \sum_j^n \sum_{k_1}^p \cdots \sum_{k_q}^p \sum_{l_1}^p \cdots \sum_{l_q}^p \left\{ \phi_{ijk_1 \dots k_q} y_{k_1} \cdots y_{k_q} \gamma_{jl_1 \dots l_q} y_{l_1} \cdots y_{l_q} \right\} = 0 \quad \forall i = 1 : m.$$

For each i , this is an algebraic form of order $q + q_g$

$$\sum_j^n \sum_{k_1 \dots k_q, l_1 \dots l_q \in \{d_1 \dots d_{q+q_g}\}} \phi_{ij d_1 \dots d_q} \gamma_{j d_{q+1} \dots d_{q+q_g}} = 0$$

$$\forall i = 1 : m, \quad k_1 = 1 : p, \quad \dots, \quad k_q = 1 : p,$$

$$l_1 = 1 : p, \quad \dots, \quad l_q = 1 : p,$$

where the second sum sums over all permutations of $k_1 \dots k_q, l_1 \dots l_q$.

D Real data experiment description

This section contains more details about the real data experiment described in Section 5.2. This data was collected together with the data previously published by Solin et al. (2015).

D.1 Experiment setup

To collect the measurements we made use of a wooden platform, see Figure 6. The platform was equipped with a Trivisio Colibri wireless IMU (TRIVISIO Prototyping GmbH, <http://www.trivisio.com/>), sampled at 100 Hz. The sensor includes both an accelerometer, a gyroscope, and a magnetometer. For additional validation a Google Nexus 5 smartphone was also mounted on the platform even though its data was never used in this experiment.

On the platform, five markers were mounted. An optical reference system (Vicon) with several cameras mounted in the ceiling measured the 3D position of each marker, and hence also the position and the orientation of the platform relative to its predefined origin.

D.2 Experiment execution

The sensor platform was moved around by hand up and down in a volume of $4 \times 4 \times 2$ meters, see Figure 5. During the experiment, measurements were collected from the sensors on the platform as well from the optical reference system. The data from the different sensors were collected asynchronously. The experiment lasted for 187 seconds.

D.3 Pre-processing of data

The position and orientation data from the optical reference system was synchronized with the data from the Trivisio sensor. The synchronization was performed based on correlation analysis of the angular velocities measured by both systems.

The position in global coordinates of the Trivisio sensor was computed based on the position data, the orientation data, and the displacement of the Trivisio sensor relative to the predefined origin of the platform.



Figure 5: Three snapshots from the measurement collection. The sensor platform was moved around by hand during approximately three minutes.

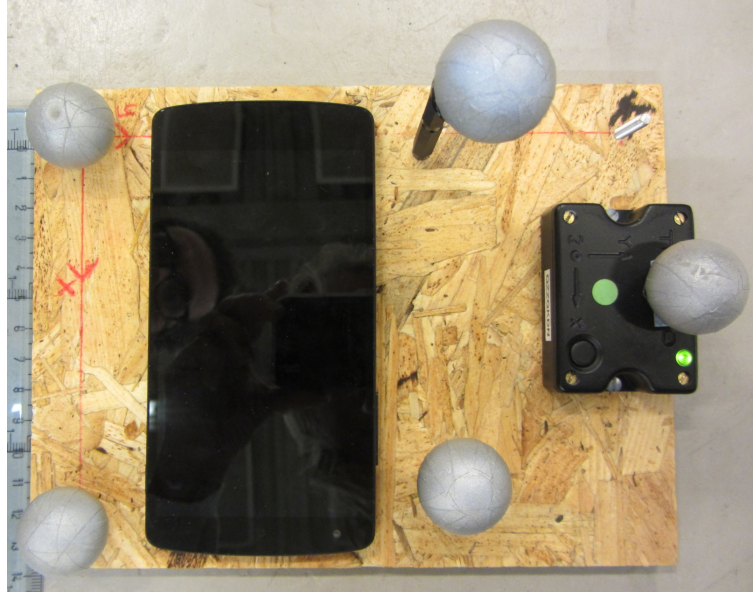


Figure 6: Platform with magnetic sensors. The sensor to the left is the Trivision sensor, whose magnetometer data we used during the experiment. The platform was also equipped with multiple markers visible to an optical reference system (Vicon).

The magnetometer data from the Trivision sensor was rotated from sensor-fixed coordinates to global coordinates using the orientation data from the optical reference system. These rotated measurements describe the magnetic field in global coordinates at the sensor positions computed above. In Section 5.2 of the main paper, these position data and magnetic field data are considered as input and output data, respectively.

References

- Abrahamsen, P. and Benth, F. E. (2001). Kriging with inequality constraints. *Math. Geol.*, 33(6):719–744.
- Álvarez, M. A., Rosasco, L., and Lawrence, N. D. (2012). Kernels for vector-valued functions: A review. *Foundations and Trends in Machine Learning*, 4(3):195–266.
- Andrade-Pacheco, R., Mubangizi, M., Quinn, J., and Lawrence, N. (2016). *Monitoring Short Term Changes of Infectious Diseases in Uganda with Gaussian Processes*, pages 95–110. Springer International Publishing.
- Constantinescu, E. M. and Anitescu, M. (2013). Physics-based covariance models for Gaussian processes with multiple outputs. *International Journal for Uncertainty Quantification*, 3(1):47–71.
- Da Veiga, S. and Marrel, A. (2012). Gaussian process modeling with inequality constraints. *Annales de la faculté des sciences de Toulouse Mathématiques*, 21(3):529–555.

- Garnett, R. (2015). Lecture 11: Bayesian quadrature. University Lecture. Accessed: 2016-12-20.
- Ghahramani, Z. (2015). Probabilistic machine learning and artificial intelligence. *Nature*, 521:452–459.
- Graepel, T. (2003). Solving noisy linear operator equations by gaussian processes: Application to ordinary and partial differential equations. In *Proceedings of the Twentieth International Conference on Machine Learning (ICML)*.
- Hennig, P. and Kiefel, M. (2013). Quasi-Newton methods: A new direction. *The Journal of Machine Learning Research*, 14(1):843–865.
- Koyejo, O., Lee, C., and Ghosh, J. (2013). Constrained Gaussian process regression for gene-disease association. *Proceedings of the IEEE 13th International Conference on Data Mining Workshops*, 00:72–79.
- Luenberger, D. G. (1969). *Optimization by vector space methods*. John Wiley & Sons, Inc.
- Maatouk, H. and Bay, X. (2016). Gaussian process emulators for computer experiments with inequality constraints. Technical report, arXiv:1606.01265.
- Navarro, A. K. W., Frellsen, J., and Turner, R. E. (2016). The multivariate generalised von mises distribution: Inference and applications. Technical report, arXiv:1602.05003.
- Nguyen, N. and Peraire, J. (2015). Gaussian functional regression for linear partial differential equations. *Computer Methods in Applied Mechanics and Engineering*, 287:69–89.
- Nguyen, N. and Peraire, J. (2016). Gaussian functional regression for output prediction: Model assimilation and experimental design. *Journal of Computational Physics*, 309:52–68.
- Papoulis, A. and Pillai, S. U. (1991). *Probability, random variables, and stochastic processes*. McGraw-Hill Education, New York.
- Rasmussen, C. E. and Williams, C. K. I. (2006). *Gaussian processes for machine learning*. MIT press, Cambridge, MA.
- Ross, J. and Dy, J. (2013). Nonparametric mixture of Gaussian processes with constraints. In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, volume 28, pages 1346–1354. JMLR Workshop and Conference Proceedings.
- Rudovic, O. and Pantic, M. (2011). Shape-constrained gaussian process regression for facial-point-based head-pose normalization. In *Proceedings of the International Conference on Computer Vision (ICCV)*.
- Salzmann, M. and Urtasun, R. (2010). Implicitly constrained Gaussian process regression for monocular non-rigid pose estimation. In *Neural Information Processing Systems (NIPS)*.
- Särkkä, S. (2011). Linear operators and stochastic partial differential equations in Gaussian process regression. In *Proceedings of the Artificial Neural Networks and Machine Learning – ICANN 2011*, pages 151–158. Springer.
- Solin, A., Kok, M., Wahlström, N., Schön, T. B., and Särkkä, S. (2015). Modeling and interpolation of the ambient magnetic field by gaussian processes. Technical report, arXiv:1509.04634.
- Tran, C., Pavlovic, V., and Kopp, R. (2015). Gaussian process for noisy inputs with ordering constraints. Technical report, arXiv:1507.00052.
- Wahlström, N. (2015). *Modeling of Magnetic Fields and Extended Objects for Localization Applications*. PhD thesis, Division of Automatic Control, Linköping University.